

1 Identifying and correcting repeat-calling errors in nanopore sequencing of 2 telomeres

3

4 Kar-Tong Tan^{1,2,3}, Michael K. Slevin^{1,4}, Matthew Meyerson^{1,2,3,4,#}, Heng Li^{5,6,#}

5 ¹Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

6 ²Cancer Program, Broad Institute of MIT and Harvard, Cambridge, MA, USA

7 ³Department of Genetics, Harvard Medical School, Boston, MA, USA

8 ⁴Center for Cancer Genomics, Dana-Farber Cancer Institute, Boston, MA, USA

9 ⁵Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA

10 ⁶Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

11

12 #Correspondence to matthew_meyerson@dfci.harvard.edu or hli@jimmy.harvard.edu

13

14

15 Abstract

16

17 Nanopore long-read genome sequencing is emerging as a potential approach for the study of
18 genomes including long repetitive elements like telomeres. Here, we report extensive
19 basecalling induced errors at telomere repeats across nanopore datasets, sequencing platforms,
20 basecallers, and basecalling models. We found that telomeres which are represented by
21 (TTAGGG)_n and (CCCTAA)_n repeats in many organisms were frequently miscalled (~40-50% of
22 reads) as (TTAAAA)_n, or as (CTTCTT)_n and (CCCTGG)_n repeats respectively in a strand-
23 specific manner during nanopore sequencing. We showed that this miscalling is likely caused by
24 the high similarity of current profiles between telomeric repeats and these repeat artefacts,
25 leading to mis-assignment of electrical current profiles during basecalling. We further
26 demonstrated that tuning of nanopore basecalling models, and selective application of the tuned
27 models to telomeric reads led to improved recovery and analysis of telomeric regions, with little
28 detected negative impact on basecalling of other genomic regions. Our study thus highlights the
29 importance of verifying nanopore basecalls in long, repetitive, and poorly defined regions of the
30 genome, and showcases how such artefacts in regions like telomeres can potentially be
31 resolved by improvements in nanopore basecalling models.

32

33

34

35 Keywords

36

37 Nanopore-sequencing, long-reads, telomere, basecalling

38

39

40 Background

41
42 Telomeres are protective caps found on chromosomal ends, and are known to play critical roles
43 in a wide range of biological processes and human diseases [1,2]. These highly repetitive
44 structures enable cells to deal with the “end-replication problem” through the action of
45 telomerase which adds telomeric repeats to the ends of chromosomes. In cancer, the
46 reactivation of telomerase to drive telomere elongation is estimated to occur in as many as 90%
47 of human cancers, and has been shown experimentally to be critical for malignant
48 transformation [3–8]. As one ages, telomeres are also known to progressively shorten, and are
49 thus thought to also play a central role in the process of aging [9–11]. In many organisms,
50 telomeres are characterized by (TTAGGG)_n repeats that vary in length of between 2 and 20kb
51 long, which are not readily resolved by short-read sequencing approaches. Given the
52 importance of telomeres in a wide range of biological process and the technical challenges
53 associated with their analysis using short-read sequencing, there is significant interest in
54 applying emerging techniques like long-read sequencing to study these repetitive structures.

55
56 Long-read sequencing has emerged as a powerful technology for the study of long repetitive
57 elements in the genome. Two main platforms, Single Molecule Real Time (SMRT) sequencing,
58 and Nanopore sequencing, have been developed to generate sequence reads of over 10
59 kilobases from DNA molecules [12,13]. In SMRT Sequencing, the incorporation of DNA
60 nucleotides is captured real time via one of four different fluorescent dyes attached to each of
61 the four DNA bases, thereby allowing the corresponding DNA sequence to be inferred.
62 Sequencing of the same DNA molecule multiple times in a circular manner further allows highly
63 accurate consensus sequence of the DNA molecule to be generated in a process termed Pacific
64 Biosciences (PacBio) High-Fidelity (HiFi) sequencing [12]. During Nanopore sequencing, the
65 ionic current, which varies according to the DNA sequence, is measured while a single-stranded
66 DNA molecule passes through a nanopore channel. The electrical current measurement is then
67 converted into the corresponding DNA sequence using a deep neural network trained on a
68 collection of ionic current profiles of known DNA sequences [13]. Notably, both platforms enable
69 long DNA molecules of more than 10 kilo-base-pairs to be routinely sequenced and are thus
70 highly suited for the study of long repetitive elements like telomeres.

71 72 73 Results and discussion

74
75 In our analysis of telomeric regions with nanopore long-read sequencing in the recently
76 sequenced and assembled CHM13 sample [14,15], we surprisingly observed that telomeric
77 regions were frequently miscalled as other types of repeats in a strand-specific manner.
78 Specifically, although human telomeres are typically represented by (TTAGGG)_n repeats
79 **(Supplementary Figure 1a)**, these regions were frequently recorded as (TTAAAA)_n repeats
80 **(Figure 1a,b, Supplementary Figure 1 and 2a)**. At the same time, when examining the reverse
81 complementary strand of the telomeres which are represented as (CCCTAA)_n repeats, we
82 instead observed frequent substitution of these regions by (CTTCTT)_n and (CCCTGG)_n repeats
83 **(Figure 1a,b, Supplementary Figure 1 and 2b,c)**. Notably, these artefacts were not observed
84 on the CHM13 reference genome [14,15], or PacBio HiFi reads from the same site **(Figure**
85 **1a,b)**, suggesting that these observed repeats are artefacts of Nanopore sequencing or the
86 base-calling process, rather than real biological variations of telomeres. Further, these repeat-
87 calling errors could be observed on all chromosomal arms for the CHM13 sample
88 **(Supplementary Figure 1b,c)**, and were thus not restricted to a single chromosomal arm. The
89 examination of each telomeric long-read also indicates that these error repeats frequently co-
90 occur with telomeric repeats at the ends of each read **(Figure 1c, Supplementary Figure 3)**.

91 Together, our results suggest that telomeric regions are frequently misrepresented as other
92 types of repeats in a strand-specific manner during Nanopore sequencing.

93
94 We then assessed if these errors are broadly observed in other studies or are specific to the
95 CHM13 dataset from the Telomere-to-Telomere consortium. To assess this, we examined the
96 previously published NA12878 and HG002 Nanopore genome sequencing datasets [12,13,16].
97 Remarkably, the same basecalling errors, TTAGGG→TTAAAA, CCCTAA→CTTCTT, and
98 CCCTAA→CCCTGG, were similarly observed at telomeres in these datasets (**Figure 1d**,
99 **Supplementary Figure 4a**), suggesting that these basecalling errors at telomeres are broadly
100 observed across multiple studies. Remarkably, between 40-60% of reads at telomeric regions in
101 these three datasets display at least one of these type of basecalling repeat artefacts for the
102 Nanopore sequencing platform (**Supplementary Figure 4b**), while these errors were not
103 observed in the PacBio HiFi datasets for the same samples (**Supplementary Figure 4b**).
104 Further, we also partitioned these datasets based on the sequencing platforms used to generate
105 them, and noted that basecalling error repeats are observed across all three nanopore
106 sequencing platforms (MinION, GridION, PromethION) (**Figure 1d**, **Supplementary Figure 4a**).
107 Together, these results show that these error repeats extend across nanopore sequencing
108 datasets and sequencing platforms.

109
110 We then questioned if these error repeats are unique to specific nanopore basecallers or
111 basecalling models. We extracted reads from chromosomal ends, and re-basecalled ionic
112 current data of these reads using different basecallers and basecalling models. Using the
113 production-ready basecaller Guppy5 (Oxford Nanopore Technologies), and the developmental-
114 phase basecaller Bonito (Oxford Nanopore Technologies), we noticed that these basecalling
115 error repeats can be readily observed across both basecallers (**Figure 1e**, **Supplementary**
116 **Figure 5 and 6**). Further, these error repeats were also observed when different basecalling
117 models were applied (**Figure 1e**). Significantly, we also observed that the “fast” basecalling
118 mode in Guppy led to almost complete loss of the (CCCTAA)_n strand (**Figure 1e**,
119 **Supplementary Figure 5a**), while the “HAC” basecalling model enabled both strands to be
120 recovered, highlighting that the basecalling model applied can affect strand-specific recovery of
121 telomeric reads. Together, these results suggest that error repeats are observable across
122 nanopore basecallers, and basecalling models.

123
124 To determine the cause for these repeat-calling errors, we then examined the ionic current
125 profiles of these repeats. We thus generated ionic current profiles of these telomeric repeats
126 and these error repeats, induced by the nanopore basecallers, using known mean current
127 values of different 6-mers (**Methods**). Remarkably, we observed a high degree of similarity
128 between current profiles between telomeric repeats and these basecalling errors (**Figure 1f**).
129 Specifically, we observed that (TTAGGG)_n telomeric repeats had a high degree of similarity with
130 the (TTAAAA)_n error repeats generated by the Bonito base-caller (Pearson correlation = 0.9928,
131 Euclidean distance=4.9934) (**Supplementary Figure 7a-c**). Similarly, (CCCTAA)_n current
132 profile also showed high similarity with (CCCTGG)_n repeats (Pearson correlation = 0.9783,
133 Euclidean distance = 4.687), and reasonably good similarity with (CTTCTT)_n repeats (Pearson
134 correlation = 0.6411, Euclidean distance = 19.384) (**Supplementary Figure 7a-c**). Together,
135 these results suggest that similarities in current profiles between repeat sequences are possible
136 causes for repeat-calling errors at telomeric repeats.

137
138 We then examined if repeat-calling errors may extend to other repetitive sequences beyond
139 telomeric sequences. To address this, we search for other repeat pairs with similar current
140 profiles that may be susceptible to these repeat-calling errors. We simulated and performed
141 pairwise comparison of current profiles for all 6-mer repeats (n=8,386,560 comparisons)

142 **(Methods)**. Using similar Pearson correlation (≥ 0.99) and Euclidean distance cutoffs (≤ 5) as
143 observed for telomeric repeat errors identified in this study (**Supplementary Figure 7a-c**), we
144 identified a further 2577 pairs of repeats with similar current profiles (**Supplementary Table 1**,
145 **Supplementary Figure 7d**). For instance, we found that $(TTAGGG)_n$ telomeric repeats also
146 showed high similarities in current profiles with repeats with single-nucleotide substitutions like
147 $(TTAAGG)_n$, $(TTAGAG)_n$ and $(TTGGGG)_n$ (**Supplementary Figure 7d,e**). Repeat sequences
148 like $(GCTGCT)_n$ and $(AACGGC)_n$ that differed drastically at the sequence level, but shared
149 similar current profiles were also observed (**Supplementary Figure 7d,f**). Further, we also
150 examined the unmappable pool of CHM13 nanopore reads after mapping it to the CHM13
151 reference assembly. Remarkably, a significant pool of reads with long $(GT)_n$ repeats were
152 readily observed (**Supplementary Figure 8**). Interestingly, $(GTGTGT)_n$ repeats were also found
153 to have high similarities in current profiles with $(CTCTCT)_n$ repeats (**Supplementary Figure 7d**,
154 **Supplementary Table 1**), suggesting that the pool of unmappable $(GT)_n$ reads may include
155 $(CT)_n$ repeats. Collectively, our results suggests that these basecalling error repeats may be
156 observed at other repetitive regions, beyond telomeres.

157
158 To resolve these basecalling errors at telomeres, we then attempted to tune the nanopore
159 basecaller by providing it with more training examples of telomeres (**Figure 2a**). Notably, model
160 training was performed with a low learning rate to ensure that the majority of the model does not
161 get affected during training while ensuring that minor adjustments in the model can be made to
162 accurately basecall telomeres. Specifically, we tuned the deep neural network model underlying
163 the Bonito basecaller by training it at a low learning rate with ground truth telomeric sequences
164 extracted from the CHM13 reference genome, and current data of the corresponding reads
165 (**Methods**). As two Nanopore PromethION runs were performed on the CHM13 dataset, we
166 used the data from one run for training (run225) and tuning of the basecaller, and held out the
167 data from the second run (run 226) for evaluation of our tuned basecaller. With this approach,
168 we see a significant improvement in the base-calls of both the telomeres, and sub-telomeric
169 regions on the training data and held out dataset with clearly observable decrease in errors on
170 the chromosomal ends (**Figure 2b, Supplementary Figure 9a-d**). Together, our results indicate
171 that a nanopore base-caller can be tuned to more accurately base-call telomeric regions by
172 providing additional training examples.

173
174 As it is computationally more efficient to redo repeat-calling only for the small fraction of
175 problematic telomeric reads rather than all reads, we developed an overall strategy to select
176 these telomeric reads for re-basecalling with the tuned Bonito+telomeres basecaller (**Figure 2c**).
177 To select telomeric reads for selective re-basecalling, we relied on an observation from the
178 CHM13 reference genome and nanopore sequencing datasets. Specifically, we noticed that
179 telomeric reads which maps to the ends of the CHM13 reference genome tend to show a high
180 frequency of telomeric, or basecalling error repeats as compared to the rest of the genome
181 (**Supplementary Figure 10**). We therefore utilized this observation to separate the non-
182 telomeric reads, from the candidate telomeric reads (**Figure 2c, Methods**). These telomeric
183 reads were then re-base-called with the tuned Bonito basecaller before being recombined with
184 the pool of non-telomeric reads. Remarkably, with this strategy, we observed a significant
185 improvement in recovery of telomeric reads with $(TTAGGG)_n$ and $(CCCTAA)_n$ repeats (from 384
186 to 476 $TTAGGG$ and 373 to 686 $CCCTAA$ reads) (**Figure 2d**). At the same time, a sharp
187 reduction of these basecalling repeat errors was also observed (151 to 17 $TTAAAA$ reads, 561
188 to 48 $CTTCTT$ reads, and 337 to 20 $CCCTGG$ reads) (**Figure 2d**). Together, these results
189 suggests that our “selective tuning” approach for fixing basecalling errors at telomeres can
190 improve recovery of telomeric reads while reducing telomeric basecalling repeat artefacts.

191

192 We further evaluated our approach for possible impact on overall basecalling accuracy. While a
193 reduction in global basecalling accuracy was observed (~1-2%) when our tuned basecaller was
194 directly applied to the full dataset, caused likely by miscalling of endogenous (CTTCTT)_n
195 genomic repeats as (CCCTAA)_n, this loss of global basecalling accuracy could be avoided by
196 applying our basecaller to telomeric reads alone. Concordant with this, we did not observe
197 changes in overall basecalling accuracy with our telomere-selective tuning approach (**Figure**
198 **2e**). These results indicate that our telomere-selective tuning approach has negligible impact on
199 basecalling accuracy for the rest of the genome.

200

201

202 **Conclusion**

203

204 In this study, we showed that basecalling errors can be widely observed at telomeric regions
205 across nanopore datasets, sequencing platforms, basecallers, and basecalling models. We
206 further showed that these strand-specific basecalling errors were likely induced by similarities in
207 current profiles between different repeat types. To resolve these basecalling errors at telomeres,
208 we devised an overall strategy to re-basecall telomeric reads using a tuned nanopore basecaller.
209 More broadly, our study highlights the importance of verifying nanopore basecalls in long,
210 repetitive and poorly defined regions of the genome. For instance, this can be done either with
211 an orthogonal platform, or at a minimum by ensuring nanopore basecalls between opposite
212 strands are concordant. In the future, we anticipate that further improvements in the nanopore
213 basecaller or basecalling model as demonstrated in this study will potentially lead to the
214 reduction or elimination of these basecalling artefacts.

215

216

217 **Methods**

218

219 ***Nanopore and PacBio Datasets***

220 Nanopore and PacBio HiFi datasets for the CHM13 sample were downloaded directly from the
221 telomere-to-telomere consortium (<https://github.com/marbl/CHM13>)

222

223 Nanopore dataset for GM12878 was obtained from the Nanopore WGS consortium
224 (<https://github.com/nanopore-wgs-consortium/NA12878/blob/master/Genome.md>). PacBio HiFi
225 dataset for GM12878 was obtained from the repository at the SRA database (SRP194450), and
226 downloaded from the following link
227 (<https://www.ebi.ac.uk/ena/browser/view/SRR9001768?show=reads>)

228

229 The HG002 PacBio HiFi and Nanopore datasets were downloaded from the Human
230 Pangenome Reference Consortium (https://github.com/human-pangenomics/HG002_Data_Freeze_v1.0). Specifically, the HG002 Data Freeze (v1.0)
231 recommended downsampled data mix was downloaded. The PacBio HiFi dataset corresponds
232 to ~34X coverage of Sequel II System with Chemistry 2.0. The Nanopore dataset corresponds
233 to 60x coverage of unsheared sequencing from 3 PromethION flow cells from Shafin et al [17].

234

235 ***Extraction of candidate telomeric reads***

236 Telomeric reads were extracted by mapping all reads to the CHM13 draft genome assembly
237 (v1.0) obtained from the telomere-to-telomere consortium using Minimap2 (version 2.17-r941).
238 Subsequent to that, reads that mapped to within 10 kilobasepairs of the start and end of each
239 autosome and X-chromosome were then extracted using SAMtools (version 1.10).

240

241 ***Co-occurrence matrix***

242 Candidate PacBio HiFi and Nanopore telomeric reads were first extracted as described above,
243 and then converted into the FASTA format using SAMtools (version 1.10). Subsequent to that,
244 custom Python scripts were used to assess if each of the reads contain at least four consecutive
245 counts of the repeat sequence of interest (e.g. (TTAGGG)₄). This information is then used to
246 generate a pair-wise correlation matrix as depicted with R in the main text.

247

248 ***Basecalling of nanopore data with different basecallers and basecalling models***

249 Basecalling of Nanopore data was done using Guppy (Version 4.4.2), Guppy (Version 5.0.16)
250 and Bonito v0.3.5 (commit d8ae5eeb834d4fa05b441dc8f034ee04cb704c69). For Guppy4, four
251 different basecalling models were applied (guppy_dna_r9.4.1_450bps_fast,
252 guppy_dna_r9.4.1_450bps_hac, guppy_dna_r9.4.1_450bps_prom_fast,
253 guppy_dna_r9.4.1_450bps_prom_hac). For Guppy 5, six different basecalling models were
254 applied (dna_r9.4.1_450bps_fast, dna_r9.4.1_450bps_hac, dna_r9.4.1_450bps_sup,
255 dna_r9.4.1_450bps_fast_prom, dna_r9.4.1_450bps_hac_prom, dna_r9.4.1_450bps_sup_prom)
256 For Bonito, the v1, v2, v3, v3.1 and default basecalling models were applied.

257

258 ***Current profiles for different repeat sequences***

259 The mean current level for different k-mers sequenced by Nanopore sequencing was obtained
260 from the k-mer models published by Oxford Nanopore
261 (https://github.com/nanoporetech/kmer_models/tree/master/r9.4_180mv_450bps_6mer).

262 Circular permutations of each 6-mer of interest was generated, and their corresponding mean
263 current level extracted from the k-mer models. The current profiles for each of the indicated
264 repeat sequences were then plotted and depicted in the figure.

265

266

267

268 ***Pairwise comparison of all possible k-mers***

269 Current profile for each 6-mer repeat sequence was generated using the published k-mer
270 models as described above. Pairwise comparisons of all possible 6-mer repeat current profiles
271 was then performed (8,386,560 pairs in total). A corresponding (i) Pearson correlation value, (ii)
272 mean-centered Euclidean distance, and (iii) mean current difference for each pair of 6-mer
273 repeat current profiles were then generated. Pairs of repeats with a Pearson correlation value \geq
274 0.99 and Euclidean distance ≤ 5 were selected as putative repeat pairs that can be miscalled.

275
276 ***Tuning of bonito model***

277 The default model from Bonito v0.3.5 (commit d8ae5eeb834d4fa05b441dc8f034ee04cb704c69)
278 was used as the base model for model tuning. The training dataset needed for the training
279 process was generated from the telomeric reads from a PromethION run in the CHM13 dataset
280 (run225). More broadly, we then generate the training dataset by matching the current profiles
281 from the Nanopore run to ground truth sequences that we extracted from the CHM13 draft
282 reference genome assembly (v1.0) using custom written code.

283
284 Specifically, these telomeric reads were first basecalled using the initial Bonito basecalling
285 model, and then mapped back to the CHM13 draft reference genome assembly (v1.0). This
286 allowed each telomeric read to be properly assigned to its corresponding chromosomal arm with
287 its sub-telomeric sequence. Nonetheless, as the telomeric region of the same read could not be
288 properly mapped to the telomeric repeats due to the repeat errors, there was difficulty in
289 assigning the nanopore current data to the correct ground truth sequences in the reference
290 genome. As such, the presume length of sequences to extract was estimated using the
291 basecalling repeat error sequences, and the same length of sequences were then extracted
292 from the CHM13 reference genome to serve as ground truth sequences. With this idea and with
293 custom Perl script, we were able to generate a set of ground truth sequences and signals for
294 model tuning. These data were then formatted into the corresponding python objects required
295 by the Bonito basecaller with custom Python scripts. Using the tune function in Bonito and with
296 our prepared training dataset, we were then able to train the basecaller to convergence.

297
298 ***Selective application of tuned basecaller to telomeric reads***

299 We applied our tuned basecaller by first extracting candidate telomeric reads for re-basecalling.
300 This was done by enumerating the total 3-mer telomeric (TTAGGG, CCCTAA) and repeat
301 artefact count (TTAAAA, CTTCTT, CCCTGG) on each read. Reads with at least 10 total counts
302 of these repeats were isolated and their readnames noted. These reads were then excluded
303 from the total pool of reads via their readnames, and basecalled separately using our tuned
304 basecaller using the fast5 data of these reads. Following basecalling with the tuned basecaller,
305 these reads were then recombined with the main pool of reads.

306
307
308
309 **Abbreviations**

310 PacBio: Pacific Biosciences
311 SMRT: Single Molecule Real Time

312
313

314 **Author information**

315

316 ***Affiliations***

317 **Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA**

318 Kar-Tong Tan, Michael K. Slevin, Matthew Meyerson

319

320 **Cancer Program, Broad Institute of MIT and Harvard, Cambridge, MA, USA**

321 Kar-Tong Tan, Matthew Meyerson

322

323 **Department of Genetics, Harvard Medical School, Boston, MA, USA**

324 Kar-Tong Tan, Matthew Meyerson

325

326 **Center for Cancer Genomics, Dana-Farber Cancer Institute, Boston, MA, USA**

327 Michael K. Slevin, Matthew Meyerson

328

329 **Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA**

330 Heng Li

331

332 **Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA**

333 Heng Li

334

335

336 ***Author contributions***

337 K.T.T. and M.S. identified issues with Nanopore sequencing of telomeres, and discovered
338 basecalling errors at telomeric regions. K.T.T. evaluated basecalling errors in Nanopore
339 sequencing datasets, and designed the overall approach for correcting basecalling errors at
340 telomeric regions with inputs from H.L. and M.M. K.T.T. wrote the initial draft of the manuscript
341 with inputs from H.L. and M.M. M.M. and H.L. jointly supervised the work. All authors read,
342 revised, and approved the submission of the manuscript.

343

344

345 ***Corresponding authors***

346 Correspondence to Matthew Meyerson (matthew_meyerson@dfci.harvard.edu) or Heng Li
347 (hli@jimmy.harvard.edu).

348 **Declarations**

349

350 ***Ethics approval and consent to participate***

351 Not applicable.

352

353 ***Consent for publication***

354 Not applicable.

355

356 ***Availability of data and materials***

357 Source code to apply and retrain the bonito basecalling model for telomeric region can be found
358 at the following link: https://github.com/ktan8/nanopore_telomere_basecall/.

359

360 The tuned bonito basecalling model can be downloaded from

361 [https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-](https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-0ab2f2848e3c/chm13_nanopore_trained_run225.zip)

362 [0ab2f2848e3c/chm13_nanopore_trained_run225.zip](https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-0ab2f2848e3c/chm13_nanopore_trained_run225.zip). A comprehensive version of

363 Supplementary Table 1 with all possible pairs of k-mers can be found at

364 [https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-](https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-0ab2f2848e3c/all_comparisions.similar_profile.txt.zip)

365 [0ab2f2848e3c/all_comparisions.similar_profile.txt.zip](https://zenodo.org/api/files/86cb9586-300f-493d-b9c4-0ab2f2848e3c/all_comparisions.similar_profile.txt.zip).

366

367

368 ***Competing interests***

369 H.L. is a consultant of Integrated DNA Technologies and on the SAB of Sentieon, Innozeen and
370 BGI. M.M. has a patent for *EGFR* mutations for lung cancer diagnosis issued, licensed, and with
371 royalties paid from LabCorp and a patent for EGFR inhibitors pending to Bayer; and was a
372 founding advisor of, consultant to, and equity holder in Foundation Medicine, shares of which
373 were sold to Roche.

374

375 ***Funding***

376 K.T.T. is supported by a PhRMA Foundation Informatics Fellowship, and a NUS Development
377 Grant from the National University of Singapore. M.M. is supported by an American Cancer
378 Society Research Professorship. This work was supported by grants from the National Human
379 Genome Research Institute (NHGRI) (Grant Nos. R01 HG010040, U01 HG010961, and U41
380 HG010972 to H.L.), and the National Cancer Institute (Grant No. R35 CA197568 to M.M.).

381

382 ***Acknowledgements***

383 We would like to thank all members of the H.L. and M.M. labs for helpful comments and
384 discussions. We would also like to thank the Telomere-to-Telomere consortium for generating
385 the CHM13 datasets used in this study.

386

387

388

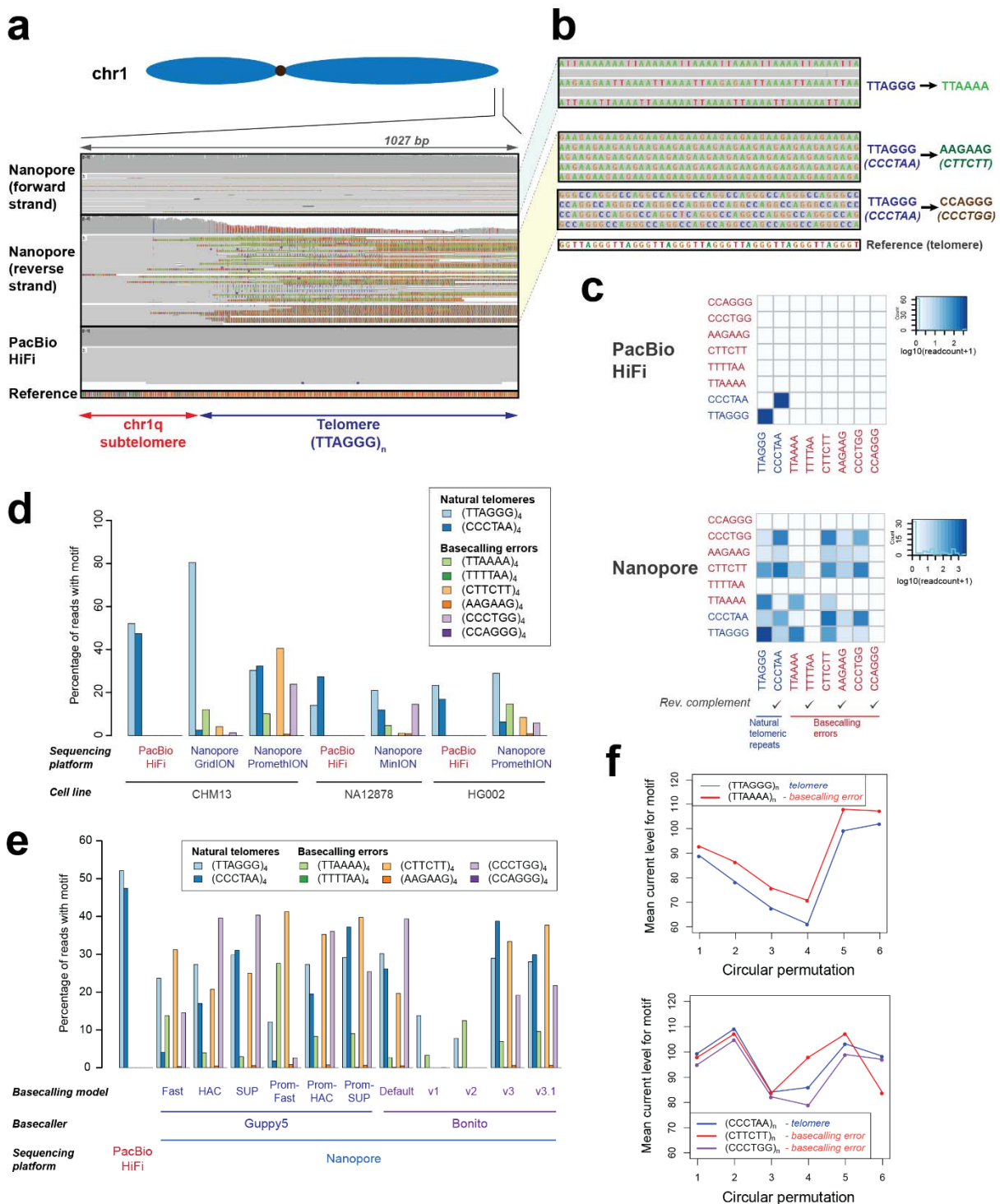
389

390 **References**

- 391 1. Shay JW, Wright WE. Telomeres and telomerase: three decades of progress. *Nat Rev Genet.*
392 2019;
- 393 2. Turner K, Vasu V, Griffin D. Telomere Biology and Human Phenotype. *Cells.* 2019;
- 394 3. Li Y, Tergaonkar V. Noncanonical functions of telomerase: Implications in telomerase-
395 targeted cancer therapies. *Cancer Res.* 2014.
- 396 4. Kim NW, Piatyszek MA, Prowse KR, Harley CB, West MD, Ho PLC, et al. Specific
397 association of human telomerase activity with immortal cells and cancer. *Science (80-).* 1994;
- 398 5. Meyerson M, Counter CM, Eaton EN, Ellisen LW, Steiner P, Caddle SD, et al. hEST2, the
399 putative human telomerase catalytic subunit gene, is up- regulated in tumor cells and during
400 immortalization. *Cell.* 1997;
- 401 6. Kolquist KA, Ellisen LW, Counter CM, Meyerson M, Tan LK, Weinberg RA, et al. Expression
402 of TERT in early premalignant lesions and a subset of cells in normal tissues. *Nat Genet.* 1998;
- 403 7. Li Y, Tergaonkar V. Telomerase reactivation in cancers: Mechanisms that govern
404 transcriptional activation of the wild-type vs. mutant TERT promoters. *Transcription.* 2016.
- 405 8. Yuan X, Larsson C, Xu D. Mechanisms underlying the activation of TERT transcription and
406 telomerase activity in human cancer: old actors and new players. *Oncogene.* 2019.
- 407 9. Shay JW. Telomeres and aging. *Curr. Opin. Cell Biol.* 2018.
- 408 10. Aubert G, Lansdorp PM. Telomeres and aging. *Physiol. Rev.* 2008.
- 409 11. Shamas MA. Telomeres, lifestyle, cancer, and aging. *Curr Opin Clin Nutr Metab Care.*
410 2011;
- 411 12. Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, et al. Accurate
412 circular consensus long-read sequencing improves variant detection and assembly of a human
413 genome. *Nat Biotechnol.* 2019;
- 414 13. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and
415 assembly of a human genome with ultra-long reads. *Nat Biotechnol.* 2018;
- 416 14. Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, et al. Telomere-to-
417 telomere assembly of a complete human X chromosome. *bioRxiv.* 2019;
- 418 15. Logsdon GA, Vollger MR, Hsieh PH, Mao Y, Liskovych MA, Koren S, et al. The structure,
419 function and evolution of a complete human chromosome 8. *Nature.* 2021;
- 420 16. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing of
421 seven human genomes to characterize benchmark reference materials. *Sci Data.* 2016;
- 422 17. Shafin K, Pesout T, Lorig-Roach R, Haukness M, Olsen HE, Bosworth C, et al. Nanopore
423 sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes.
424 *Nat Biotechnol.* 2020;38.

425
426
427

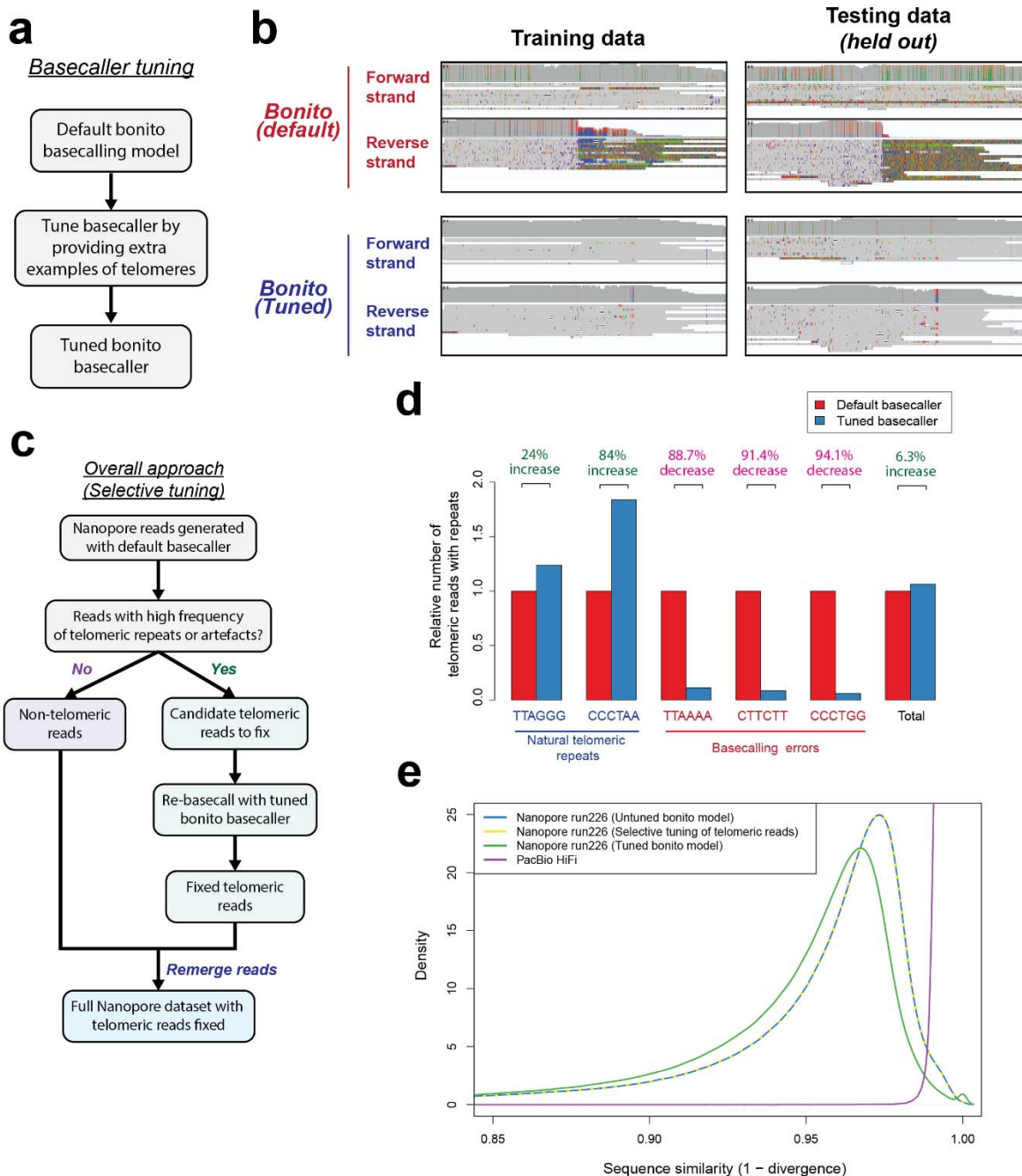
Figure 1



428
 429 **Figure 1 Strand-specific Nanopore basecalling errors are pervasive at telomeres. (a,b)**
 430 IGV screenshot illustrating the three types of basecalling errors found on the forward and
 431 reverse strands of telomeres for Nanopore sequencing. (TTAGGG)_n on the forward strand of
 432 Nanopore sequencing data was basecalled as (TTAAAA)_n while (CCCTAA)_n on the reverse

433 strand was basecalled as $(CTTCTT)_n$ and $(CCCTGG)_n$. PacBio HiFi data generated from the
434 same cell line (CHM13) is depicted as a control. Reference genome indicated in the plot
435 corresponds to the chm13 draft genome assembly (v1.0). **(c)** Co-occurrence heatmap
436 illustrating the frequency of co-occurrence of repeats corresponding to natural telomeres, or to
437 basecalling errors in PacBio HiFi and Nanopore long-reads found at chromosomal ends (within
438 10kb of annotated end of the reference genome). Diagonal of co-occurrence matrix represents
439 counts of long-reads with only a single type of repeats observed. **(d)** Basecalling errors at
440 telomeres are observed across different nanopore datasets and sequencing platforms. **(e)**
441 Basecalling errors at telomeres are observed different nanopore basecallers and basecalling
442 models. Guppy5 and the Bonito basecallers, and different basecalling models for each basecaller,
443 were used to basecall telomeric reads in the CHM13 PromethION dataset (reads that mapped
444 to flanking 10kb regions of the CHM13 reference genome). **(f)** Basecalling errors share similar
445 nanopore current profiles as telomeric repeats. Current profiles for telomeric and basecalling
446 error repeats were plotted based on known mean current profiles for each k-mer (Methods).
447
448

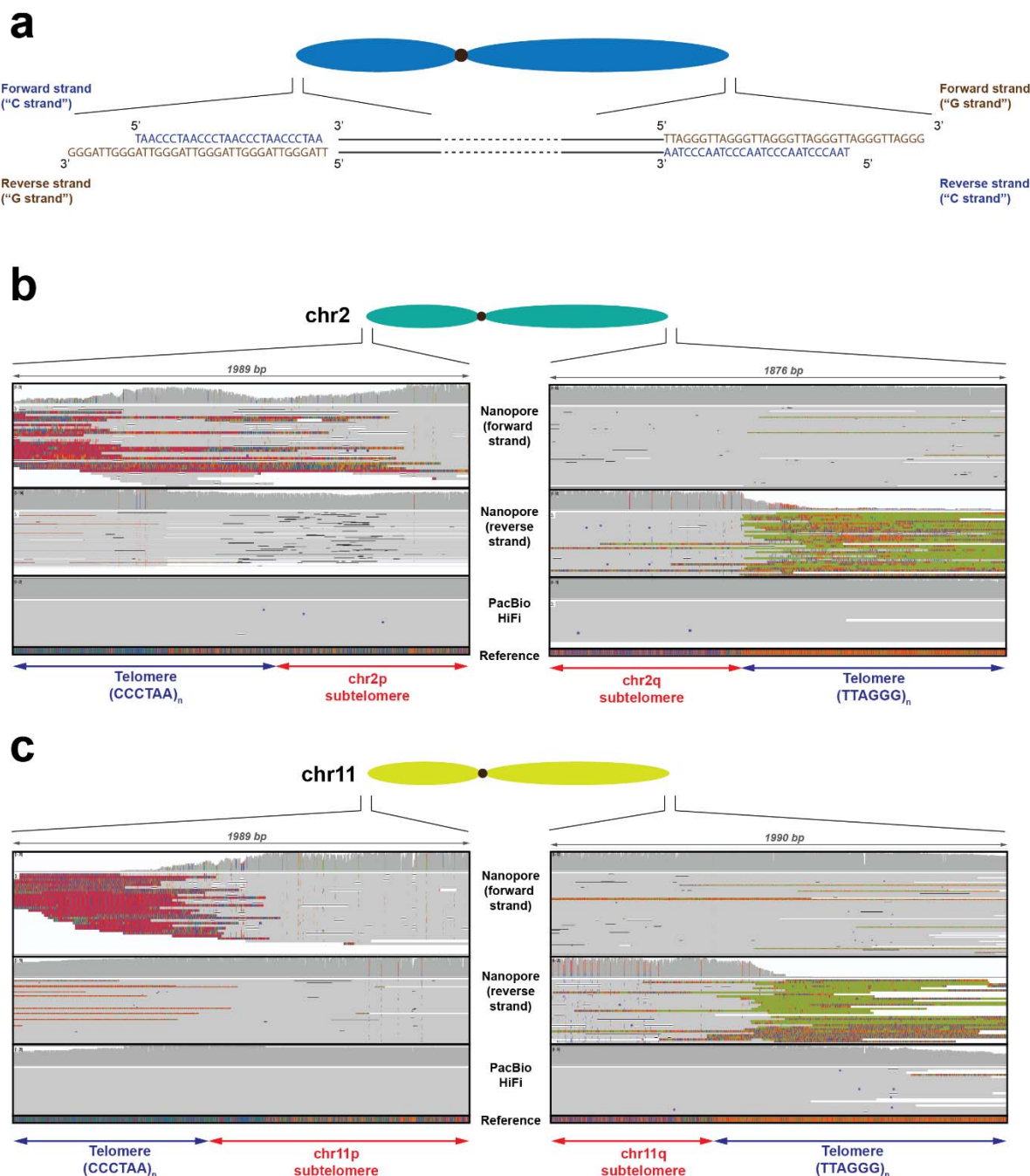
Figure 2



449
 450 **Figure 2 Selective re-basecalling of telomeric reads resolves basecalling errors at**
 451 **telomeres.** (a) Approach for tuning the bonito basecalling model for improving basecalls at
 452 telomeres. (b) Tuned bonito basecalling model leads to improvement in basecalls at telomeric
 453 regions. IGV screenshots of telomeric region (chr2q) in the CHM13 dataset basecalled using the
 454 default bonito basecaller, and the tuned bonito basecalling model is as depicted. (c) Overall
 455 approach for selecting and fixing telomeric reads in nanopore sequencing datasets. Telomeric
 456 reads are selected (Methods), and rebasecalled using the tuned bonito basecalling model. (d)
 457 The selective tuning approach leads to improved recovery of telomeric reads, and decrease in

458 the number of reads with basecalling artefacts. Evaluation was performed on the held out test
459 dataset (run226). **(e)** The 'selective basecalling' approach leads to little detected negative
460 impact on basecalling of other genomic regions. The sequence similarity of all reads to the
461 reference genome for three approaches for basecalling of nanopore reads was evaluated. They
462 are applying the default bonito basecalling model to all reads (untuned bonito model), applying
463 the tuned bonito basecalling model to all reads (tuned bonito model), and applying the tuned
464 bonito basecalling model selectively to telomeric reads only (selective tuning of telomeric reads).
465 The density plot depicts the sequence similarity of each read against the CHM13 reference
466 genome as assessed using minimap2.
467

Supplementary Figure 1



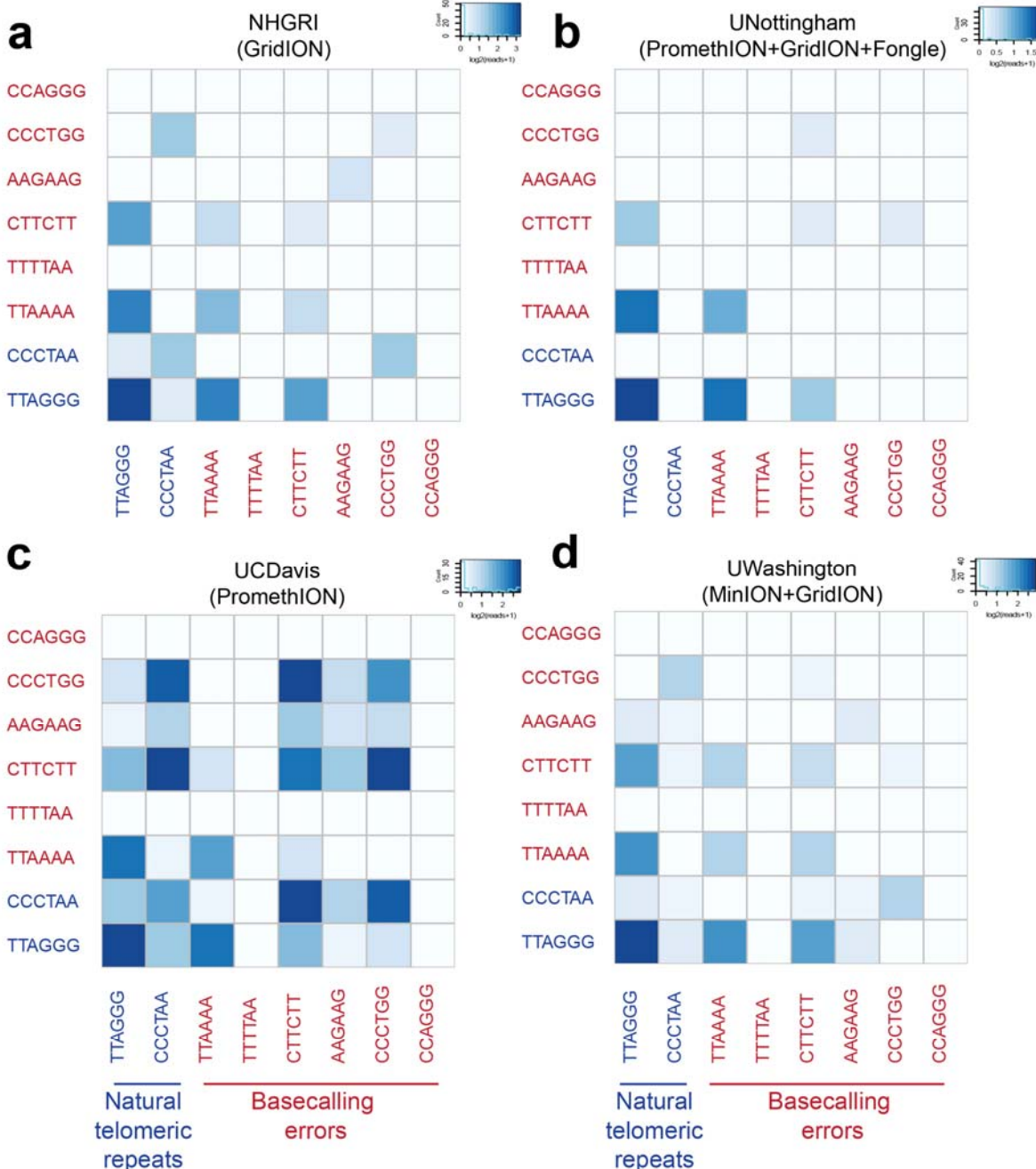
468
 469
 470
 471
 472
 473
 474
 475
 476

Supplementary Figure 1 Additional screenshots of basecalling repeat errors found on different chromosomal arms. (a) Schematic depicting sequence and orientation of telomeric repeat sequences on the p-arms (arm on the left in the schematic) and q-arms (arm on the right of the schematic) of a chromosome. Note that the forward strand for the arm on the left, and reverse strand for the arm on the right are "C-rich strands" and characterized by (CCCTAA)_n repeats in a 5'-to-3' direction. Also note that the reverse strand for the arm on the left, and forward strand for the arm on the right are "G-rich strands" and characterized (TTAGGG)_n

477 repeats in a 5'-to-3' direction. **(b-c)** Screenshots depicting additional representative examples of
478 chromosomal arms with basecalling error repeats. These are **(b)** chromosome 2 and **(c)**
479 chromosome 11. Screenshots were extracted from the Integrative Genomics Viewer for the
480 CHM13 long-read dataset mapped against the CHM13 reference genome. Related to Figure 1a.
481

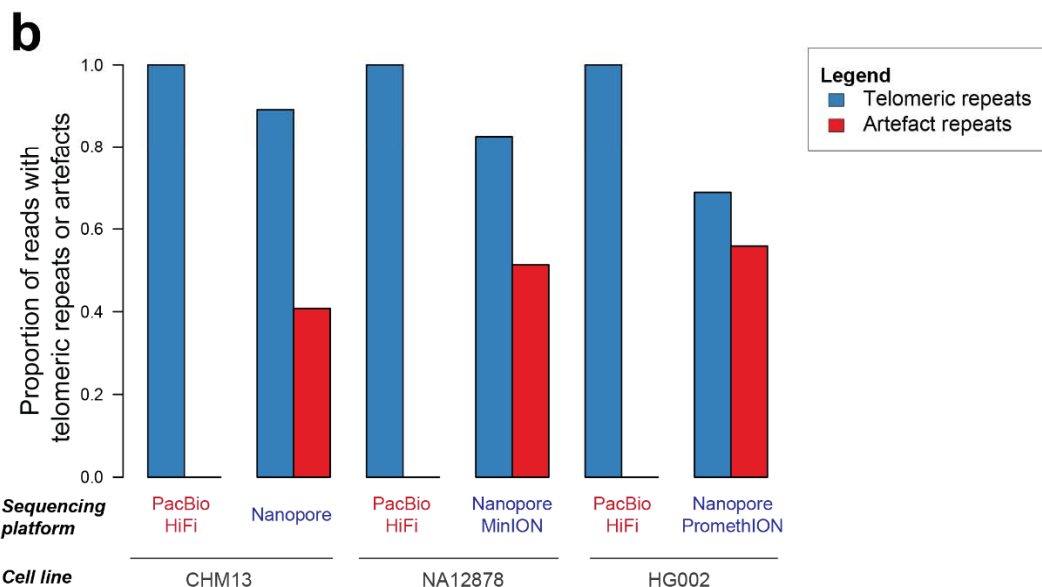
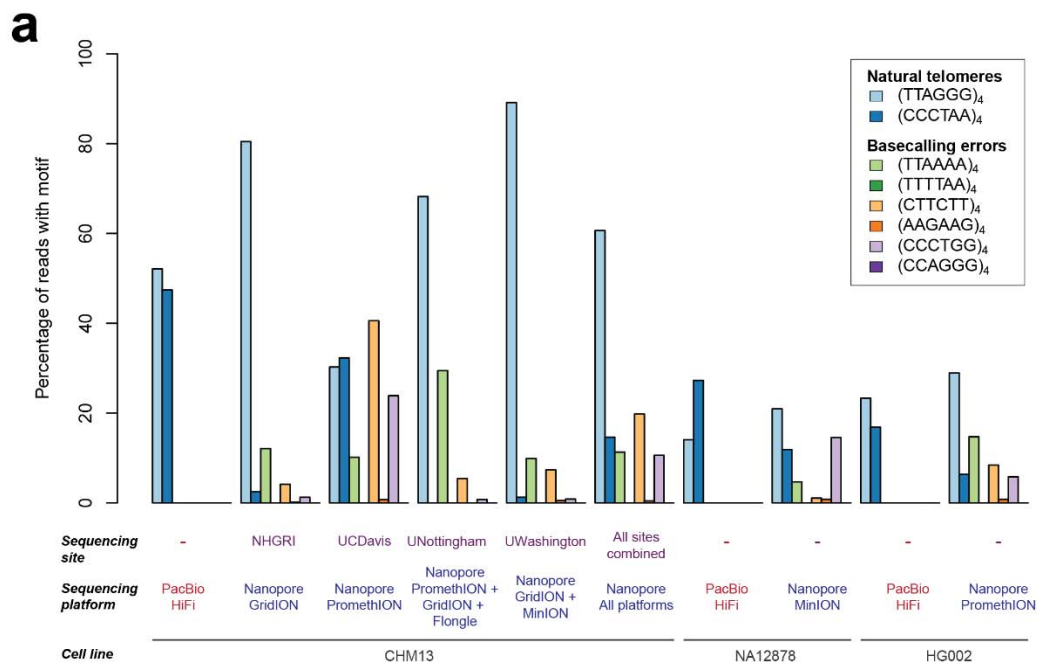
488 (CCCTAA)_n to (CCCTGG)_n. Note that **(b)** and **(c)** represents the reverse complementary
489 sequence of the actual nanopore long-read sequence. Also note that the repeats were found on
490 the end of each read as expected given that telomeric repeats are typically found on the end of
491 the chromosomes.
492

Supplementary Figure 3



493
 494 **Supplementary Figure 3 Co-occurrence heatmap illustrating the frequency of co-**
 495 **occurrence of telomeric repeats and basecalling errors for the CHM13 Nanopore dataset**
 496 **generated at different sites.** These are **(a)** National Human Genome Research Institute
 497 **(NHGRI), (b)** University of Nottingham (UNottingham), **(c)** University of California, Davis
 498 **(UCDavis) and (d)** University of Washington (UWashington). The sequencing platforms used for
 499 sequencing at each of the sites are also as indicated. This figure is related to Figure 1b.
 500

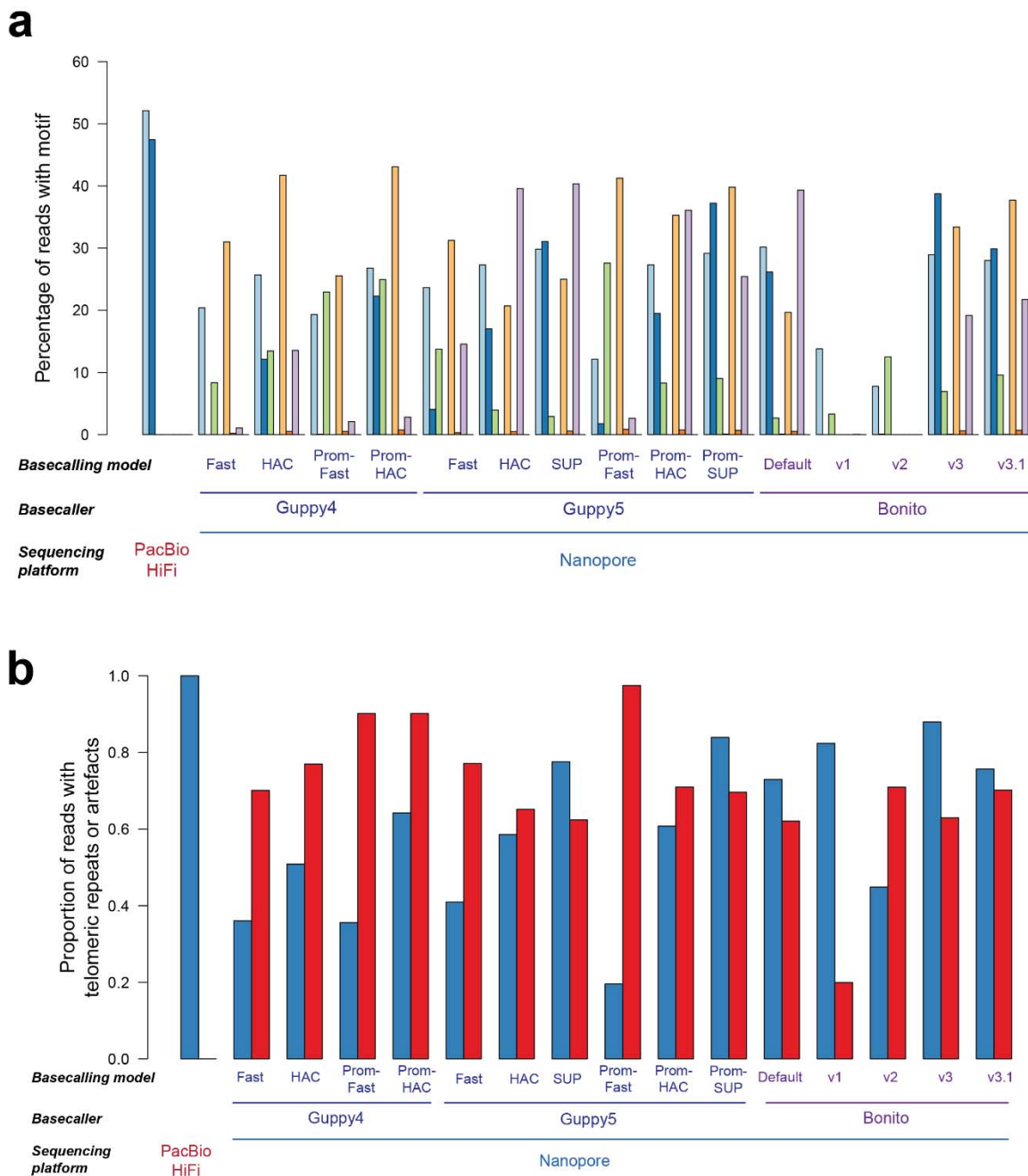
Supplementary Figure 4



501
 502 **Supplementary Figure 4 Frequency of telomeric repeat errors in different Nanopore**
 503 **sequencing dataset and sequencing platforms. (a)** Frequency of basecalling error repeats
 504 on three different cell lines generated by different Nanopore sequencing platforms. This figure
 505 is an extension Figure 1d. **(b)** Aggregated fraction of basecalling error repeats for different cell
 506 lines and datasets.

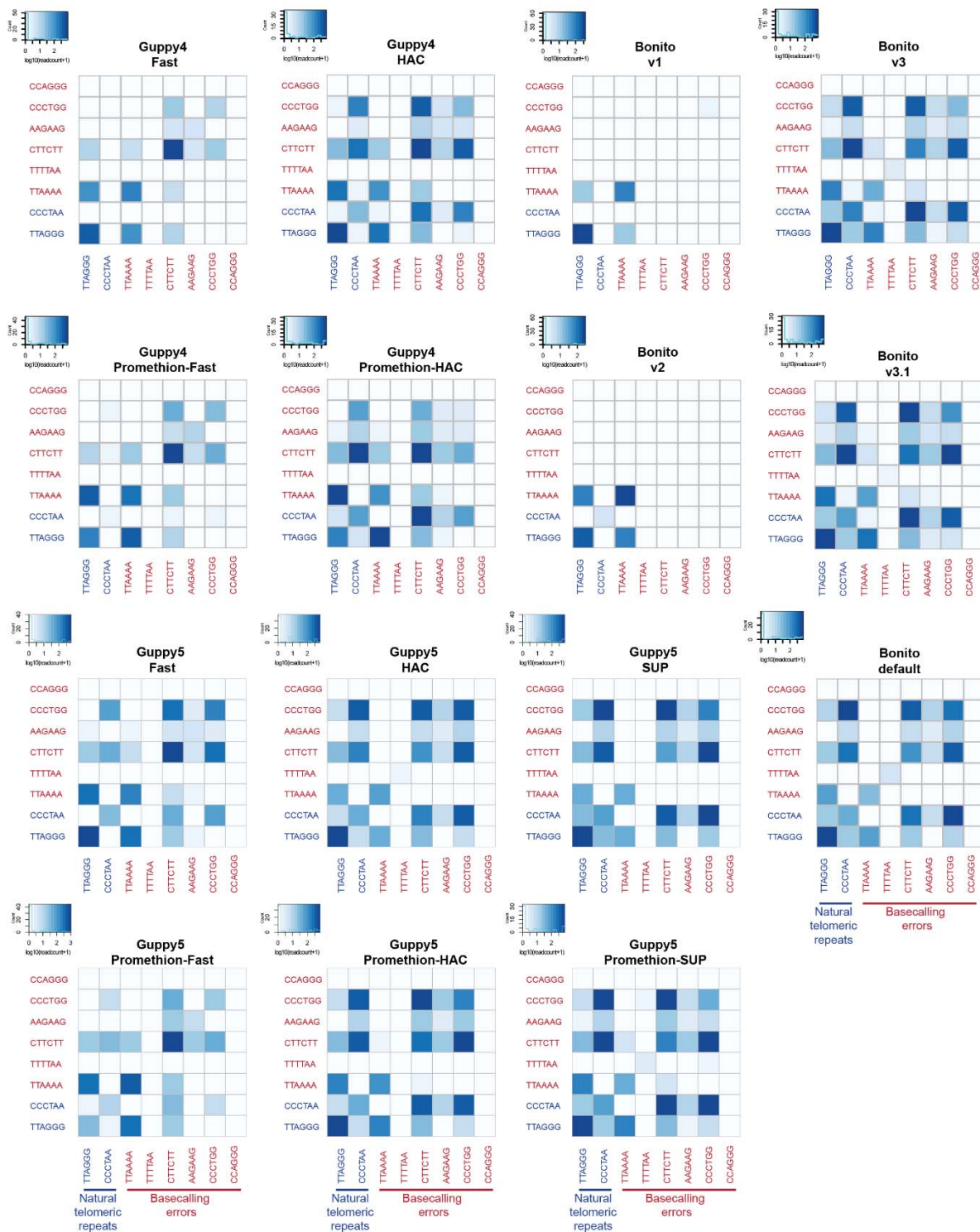
507
 508
 509

Supplementary Figure 5



510
 511 **Supplementary Figure 5 Frequency of telomeric repeat errors in different Nanopore**
 512 **basecallers. (a)** Frequency of basecalling error repeats for different basecallers (Guppy4,
 513 Guppy5 and Bonito) and basecalling models. This figure is an extension of Figure 1e. **(b)**
 514 Aggregated fraction of basecalling error repeats for different basecallers and basecalling models.
 515
 516

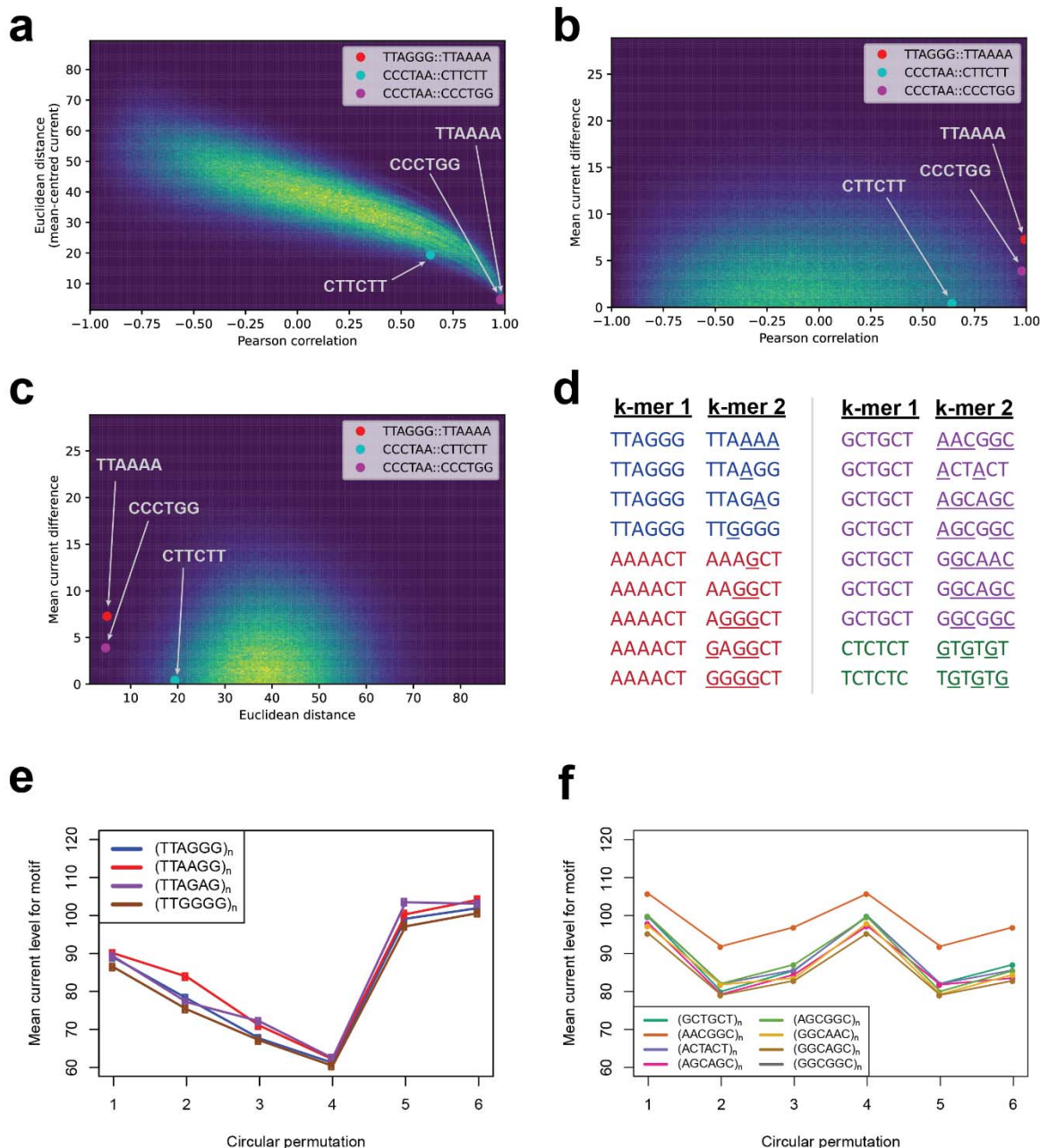
Supplementary Figure 6



517
518
519
520
521

Supplementary Figure 6 Co-occurrence heatmap for different Nanopore basecalling models. Different nanopore basecallers and basecalling models were applied to the CHM13 Nanopore promethion datasets. The frequency of telomeric repeats and basecalling artefacts observed on reads obtained are as depicted.

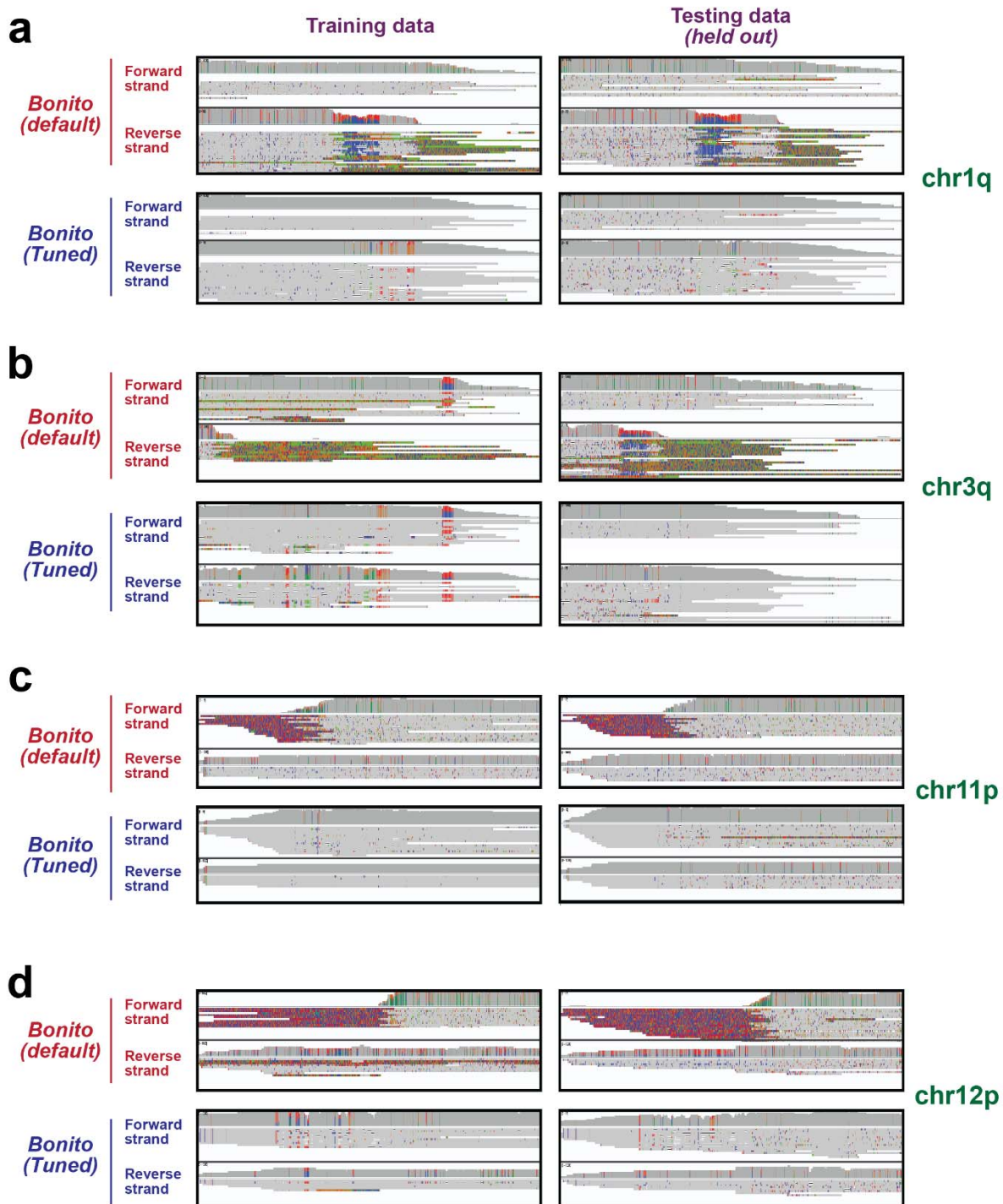
Supplementary Figure 7



522
523 **Supplementary Figure 7 Similarities between current profiles for all possible pairs of 6-**
524 **mer repeats. (a-c)** Heatmaps depicting the Euclidean distances, Pearson correlation, and mean
525 current differences between current profiles between all possible 6-mer repeat sequences.
526 These are depicted as pairwise plots for **(a)** the Euclidean distances vs. the Pearson correlation,
527 **(b)** the mean current difference vs. the Pearson correlation, and **(c)** the mean current difference
528 vs. the Euclidean distance. The pairwise comparisons between the telomeric repeats and the
529 observed basecalling repeat artifacts are also highlighted in the plots. **(d)** Example pairs of k-
530 mer repeats with similar current profiles are as indicated. The nucleotides in k-mer 2 that differs
531 from k-mer 1 is underlined to highlight the nucleotides that differ between the two types of

532 repeats. **(e-f)** Current profiles for repeats which were predicted to be highly similar to each other.
533 These are depicted for **(e)** TTAGGG telomeric repeats and telomere-like repeat sequences and
534 **(f)** GCTGCT repeat sequences that were highlighted in purple in Supplementary Figure 7d.
535
536

Supplementary Figure 9

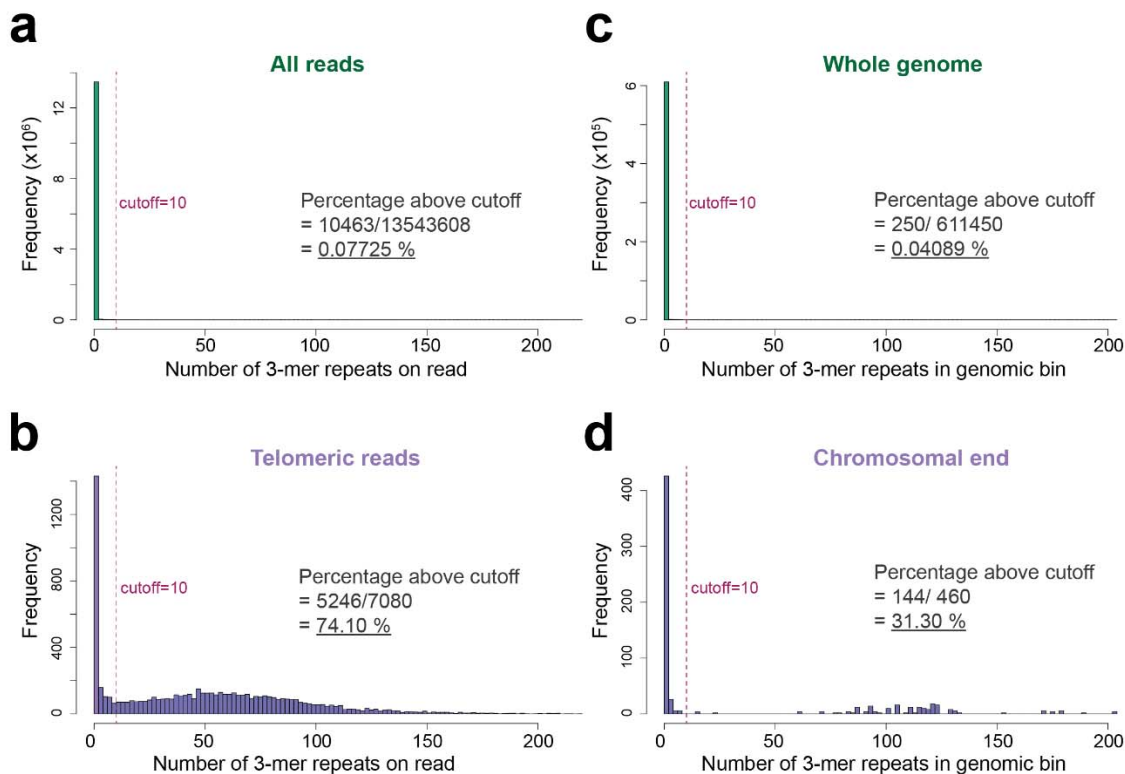


543
544
545
546
547

Supplementary Figure 9 Additional examples for the performance of the tuned bonito basecaller on telomeres on other chromosomal arms. The tuned model was applied to the training dataset used for model training, and on an additional held out test dataset that was not used during model training. IGV screenshots of the default and tuned bonito basecaller on the

548 training and testing dataset for the chromosomal arms **(a)** chr1q, **(b)** chr3q, **(c)** chr11p and **(d)**
549 chr12p are as depicted. Related to Figure 2b.
550
551

Supplementary Figure 10



552
553 **Supplementary Figure 10 Histograms depicting the frequencies of 3-mer repeats on**
554 **reads at telomeres and on reads found at the rest of the genome in the CHM13 dataset.**
555 **(a-b)** The sum of 3-mer telomeric repeats [(TTAGGG)₃ (CCCTAA)₃] and basecalling error
556 repeats [(TTAAAA)₃, (TTTTAA)₃, (CTTCTT)₃, (AAGAAG)₃, (CCCTGG)₃, (CCAGGG)₃] on **(a-b)**
557 each long-read or **(c-d)** genomic bin are as depicted on the x-axis of each histogram. The
558 histograms represent the frequency of these repeats on **(a)** all long-reads in the CHM13 dataset,
559 **(b)** telomeric reads in the CHM13 dataset, **(c)** 20 kb genomic bins with 10 kb moving window for
560 the full CHM13 reference genome, **(d)** and for the 10 genomics bins on each chromosomal end
561 of the CHM13 genome.

562
563

564 **Supplementary Tables**

565

566 **Supplementary Table 1 List of k-mers with high similarities in current profiles.** The
567 pearson correlation, Euclidean distance, and mean current difference between each pair of k-
568 mer is as presented in the table.